

SPEECH RECOGNITION DEVICE AND SPEECH RECOGNITION METHOD

Patent Number: JP2001255886
Publication date: 2001-09-21
Inventor(s): SUZUKI SATORU; OONO TAKEO; KIMURA TATSUYA
Applicant(s): MATSUSHITA ELECTRIC IND CO LTD
Requested Patent: JP2001255886
Application Number: JP20000064919 20000309
Priority Number(s):
IPC Classification: G10L15/10; G10L11/06; G10L15/02; G10L15/20; G10L21/02
EC Classification:
Equivalents:

Abstract

PROBLEM TO BE SOLVED: To provide a speech recognition device by which a user does not utter specific contents is not and which uses a speaker normalization processing capable of speedily performing normalization into individual features of a speaker without requiring an on-line 'non-teacher'.
SOLUTION: Featured values such as LPC cepstrum coefficients are extracted with the speech digitalized by A/D conversion as input signals (S10). Then, frequency axis conversion is conducted for the featured value such as the LPC cepstrum in order to normalize the effect caused by the individuality of the length of the vocal track of the uttering person (S30). Then, a matching is conducted between the featured values of the inputted speech that is frequency axis converted and the acoustic model featured values beforehand learned from plural speakers (S50). After that, the inputted utterings are made as teacher's signals based on the recognition result computed in the S50 and optimum conversion coefficients are obtained (S60). Then, smoothing is conducted for the conversion coefficients to absorb dispersion caused by the speakers and the phonemes and new updated frequency axis conversion coefficients are obtained (S70).

Data supplied from the esp@cenet database - I2

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2001-255886

(P2001-255886A)

(43) 公開日 平成13年9月21日 (2001.9.21)

(51) Int.Cl. ⁷	識別記号	F I	テーマコード* (参考)
G 1 0 L 15/10		G 1 0 L 101: 16	5 D 0 1 5
11/06		9/16	3 0 1 B 9 A 0 0 1
15/02		3/00	5 1 5 D
15/20		3/02	3 0 1 A
21/02			

審査請求 未請求 請求項の数 9 O L (全 13 頁) 最終頁に続く

(21) 出願番号 特願2000-64919 (P2000-64919)

(22) 出願日 平成12年3月9日 (2000.3.9)

(71) 出願人 000005821

松下電器産業株式会社

大阪府門真市大字門真1006番地

(72) 発明者 鈴木 哲

神奈川県川崎市多摩区東三田3丁目10番1号 松下技研株式会社内

(72) 発明者 大野 剛男

神奈川県川崎市多摩区東三田3丁目10番1号 松下技研株式会社内

(74) 代理人 100097445

弁理士 岩橋 文雄 (外2名)

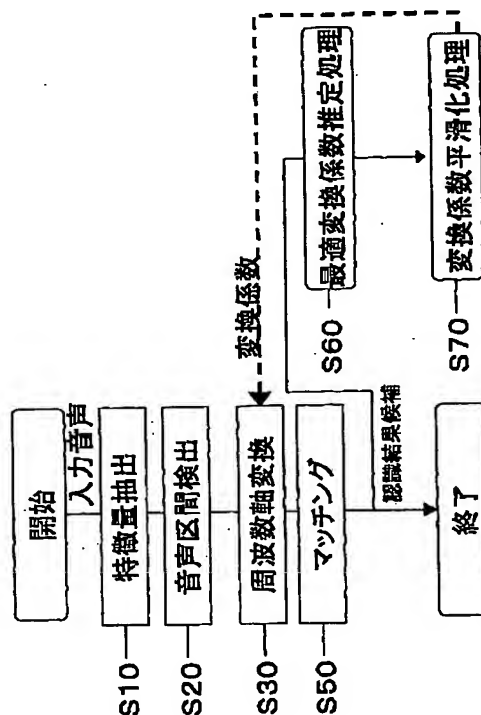
最終頁に続く

(54) 【発明の名称】 音声認識方法および音声認識装置

(57) 【要約】

【課題】 利用者に対して一定内容の発声を促す必要性が無い、オンライン「教師なし」で話者の個人性特徴にすばやく正規化できる話者正規化処理を用いる音声認識装置を提供する。

【解決手段】 A/D変換を行ってデジタル化された音声を入力信号として、LPCケプストラム係数等の特徴量を抽出し (S10)、発声者の声道長の個人性に起因する影響を正規化するために、LPCケプストラム等の特徴量に周波数軸の変換を施し (S30)、周波数軸変換を施された入力音声の特徴量と予め複数話者から学習した音響モデル特徴量とのマッチングを行なう (S50)。その後、S50において算出された認識結果をもとに入力発声を教師信号として最適な変換係数を求め (S60)、話者や音韻によるばらつきを吸収するため変換係数平滑化を行い、新たな周波数軸変換係数を更新する (S70)。



【特許請求の範囲】

【請求項1】 入力音声の特徴量を抽出する特徴量抽出ステップと、前記入力音声の特徴量の周波数軸を少なくとも1つの周波数軸変換係数から構成される周波数軸変換係数列を用いて変換する周波数軸変換ステップと、前記周波数軸変換を施した入力音声の特徴量と予め複数の話者から学習した音響モデル特徴量とをマッチングし、認識結果候補を出力するマッチングステップと、前記認識結果候補のうち少なくとも1つから表現される音素系列に対して少なくとも1つの周波数軸変換係数から構成される最適な周波数軸変換係数列を推定する最適変換係数推定ステップと、前記求められた最適な周波数軸変換係数列と保持された過去に求められた周波数軸変換係数列とを平滑化し、新たな周波数軸変換係数列を更新・保持する変換係数平滑化ステップとを有することを特徴とする音声認識方法。

【請求項2】 変換係数平滑化ステップは、少なくとも1つの周波数軸変換係数から構成される最新の周波数軸変換係数列と、保持された過去に求められた少なくとも1つの周波数軸変換係数から構成される周波数軸変換係数列とを比較することによって話者の交代を検知することを特徴とする請求項1記載の音声認識方法。

【請求項3】 入力音声の特徴量を抽出する特徴量抽出ステップと、前記入力音声から無声音／有声音区間を弁別検出する音声区間検出ステップと、前記入力音声の特徴量の周波数軸を前記無声音／有声音区間情報に応じて周波数軸変換係数列を用いて変換する周波数軸変換ステップと、前記周波数軸変換を施した入力音声の特徴量と予め複数の話者から学習した音響モデル特徴量とをマッチングし、認識結果候補を出力するマッチングステップと、前記認識結果候補のうち少なくとも1つから表現される音素系列に対して少なくとも1つの周波数軸変換係数から構成される最適な周波数軸変換係数列を推定する最適変換係数推定ステップと、前記求められた最適な周波数軸変換係数列と保持された過去に求められた周波数軸変換係数列とを平滑化し、新たな周波数軸変換係数列を更新・保持する変換係数平滑化ステップとを有することを特徴とする音声認識方法。

【請求項4】 最適変換係数推定ステップは、前記認識結果候補を求めた際に使用した少なくとも1つの周波数軸変換係数から構成される周波数軸変換係数列を基に複数の周波数軸変換係数列候補を設定し、それぞれの周波数軸変換係数毎に、前記マッチングステップの認識結果候補のうち少なくとも1つから表現される音素系列に対して、周波数軸変換を施して得られた入力音声特徴量と予め複数の話者から学習した音響モデル特徴量とマッチングにより尤度を求め、求めた尤度のうちで最大尤度を与える少なくとも1つの周波数軸変換係数から構成される周波数軸変換係数列を選択することを特徴とする請求項1から3のいずれかに記載の音声認識方法。

【請求項5】 変換係数平滑化ステップは、未知話者に最適な周波数軸変換係数を平滑化する際に、保持された過去に求められた周波数軸変換係数の平均値と最新の周波数軸変換係数との距離を求め、所定の距離を満たす場合にのみ最新の係数として採用し、新たな周波数軸変換係数として出力することを特徴とする請求項1から4のいずれかに記載の音声認識方法。

【請求項6】 変換係数平滑化ステップは、未知話者に最適な周波数軸変換係数を平滑化する際に、保持された過去に求められた周波数軸変換係数の平均値と最新の周波数軸変換係数との距離を求め、所定の距離を満たさない場合には、過去の周波数軸変換係数を任意の初期値を与えることによって初期化を行うことを特徴とする請求項2記載の音声認識方法。

【請求項7】 入力音声の特徴量を抽出する特徴量抽出手段と、前記入力音声の特徴量の周波数軸を周波数軸変換係数列を用いて変換する周波数軸変換手段と、前記周波数軸変換を施した入力音声の特徴量と予め複数の話者から学習した音響モデル特徴量とをマッチングし、認識結果候補を出力するマッチング手段と、前記認識結果候補のうち少なくとも1つから表現される音素系列に対して少なくとも1つの周波数軸変換係数から構成される最適な周波数軸変換係数列を推定する最適変換係数推定手段と、前記求められた最適な周波数軸変換係数列と保持された過去に求められた周波数軸変換係数列とを平滑化し、新たな周波数軸変換係数列を更新・保持する変換係数平滑化手段とを有することを特徴とする音声認識装置。

【請求項8】 変換係数平滑化手段は、少なくとも1つの周波数軸変換係数から構成される最新の周波数軸変換係数列と、保持された過去に求められた少なくとも1つの周波数軸変換係数から構成される周波数軸変換係数列とを比較することによって話者の交代を検知することを特徴とする請求項7記載の音声認識装置。

【請求項9】 入力音声の特徴量を抽出する特徴量抽出手段と、前記入力音声から無声音／有声音区間を弁別検出する音声区間検出手段と、前記入力音声の特徴量の周波数軸を前記無声音／有声音区間情報に応じて周波数軸変換係数列を用いて変換する周波数軸変換手段と、前記周波数軸変換を施した入力音声の特徴量と予め複数の話者から学習した音響モデル特徴量とをマッチングし、認識結果候補を出力するマッチング手段と、前記認識結果候補のうち少なくとも1つから表現される音素系列に対して少なくとも1つの周波数軸変換係数から構成される最適な周波数軸変換係数列を推定する最適変換係数推定手段と、前記求められた最適な周波数軸変換係数列と保持された過去に求められた周波数軸変換係数列とを平滑化し、新たな周波数軸変換係数列を更新・保持する変換係数平滑化手段とを有することを特徴とする音声認識装置。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、不特定話者の音声を生認識させる分野等に利用される音声認識方法および音声認識装置に関する。

【0002】

【従来の技術】人の音声の音響的特徴は、主に音を発生させる喉つまり音源と、その音が反響しながら伝播する声道およびその形状とで構成される発声器官によって決定される。つまり話者の音響的特徴の違いは、これら話者の発声器官の特徴が主に起因して生じていると考える事ができる。

【0003】そのため、特に不特定話者を対象とした音声認識装置において、音響モデルで表現される話者集団の発声器官の特徴と比較して特異な特徴を持つ話者の認識率が低下することがあると指摘されていた。

【0004】そこで、この発声器官の特徴つまり個人性を要因とした認識率の低下を防ぎ、かつ高い認識率を保持することを目的として、話者適応化手法あるいは話者正規化手法が提案されてきた。

【0005】従来提案されてきた話者適応化、話者正規化手法としては、音響モデルパラメータを既知の音響モデルパラメータを用いて再評価することにより音響モデル自身を話者にあわせて更新あるいは選択する話者適応方法と、個々の話者の特徴空間を変換して音響モデル学習話者から表現される特徴空間にマッピングする話者正規化方法との2つに大別できる。

【0006】前者は、たとえば電子通信情報学会SP92-16(1992年)に紹介されたベクトル場平滑化法のように、適応化音声の量が増すに従い、話者適応システム性能は話者依存での学習時の性能に近づくという特徴を持つため格段の認識性能の向上が期待できるものの、適応の効果が現れるのに十分な学習音声量を獲得するまでに時間を要するという欠点がある。

【0007】後者は、たとえば論文「Frequency Warpin gによる話者の正規化、松本、脇田、日本音響学会音声研究会資料S79-25,1979-7」においては、周波数正規化スペクトルによって声道長正規化に効果があると主張しており、さらに、論文「LPC距離尺度における周波数正規化に関する検討、小林、松本、日本音響学会講演論文集1-1-5、昭和58年10月」においては、LPCスペクトルの周波数軸を伸縮する方法として一次の全域通過フィルターを用いた方法を提案している。

【0008】この後者の方法は、変換係数の変更により話者の個人性を正規化できるという特徴を持つことから、オンラインの話者適応化・正規化方法としては前者に比べて、必要とする音声量がより少ない点で実用上有効であると考えられる。さらに、話者正規化手法として、特開平6-214596号公報において、声帯音源特性に関する音声スペクトル傾斜の変動と、声道特性

(例えば声道長)に関する音声スペクトルの周波数軸方向の伸縮という個人性を同時に正規化する方法が提案されている。

【0009】以下、この従来例の音声認識装置について図9を用いて説明する。

【0010】図9の音声認識装置は、入力された音声信号の周波数特性を補正する周波数特性補正部10と、入力音声信号のケプストラム係数を入力音声特徴量として抽出する特徴量抽出部20と、入力音声信号に対し周波数軸の変換を施す周波数軸変換部30と、入力された音声信号の区間を検出する音声区間検出部40と、標準音声信号の特徴量が標準音声特徴量として予め記憶されている標準音声記憶部50と、入力音声信号に対し周波数特性補正部10、特徴量抽出部20、周波数軸変換部30により得られた入力音声特徴量と標準音声記憶部50に記憶されている標準音声特徴量との照合(マッチング)を行なうマッチング部60とから構成されている。

【0011】ところで、この音声認識装置では、不特定話者の音声をも良好に認識させることを目的として、実際の音声認識処理とその音声認識処理の開始に先立って話者適応学習処理がなされる。この2種類の処理を1つの装置で行なわせるため、図9の装置には、この装置の動作、機能を話者適応フェーズと音声認識フェーズとのいずれかに切替えるためのフェーズ選択部90がさらに設けられている。

【0012】また、これと関連させて、標準音声記憶部50には、話者適応処理用の標準音声特徴量と音声認識用の標準音声特徴量とが記憶されている。また、周波数特性補正部10には、話者適応学習用に、互いに異なる複数の周波数特性補正係数が予め用意され、また、周波数軸変換部30には、話者適応学習用に、互いに異なる複数の周波数軸変換係数が用意されている。

【0013】次に、話者適応フェーズについて説明する。話者適応フェーズにおいては、未知話者に既知の発声内容を発生させるようになっており、周波数特性補正部10、周波数軸変換部30では、この音声信号に対して、各々、複数の周波数特性補正係数、複数の周波数軸変換係数を順次に変えて処理を行ない、マッチング部60は、それぞれの場合について、周波数特性補正部10、特徴量抽出部20、周波数軸変換部30により得られた入力音声特徴量を標準音声記憶部50に記憶されている話者適応処理用の標準音声特徴量とマッチングして、各入力音声特徴量と標準音声特徴量との尤度を求め、そのうち最大尤度を与える周波数特性補正係数と周波数軸変換係数とを選択し決定するようになっている。

【0014】次に、音声認識フェーズについて説明する。音声認識フェーズにおいては、未知話者(実際には、話者適応フェーズで入力を行なった話者)の未知の発声内容の音声信号に対して、周波数特性補正部10、周波数軸変換部30では、上記話者適応フェーズにおい

て選択、決定された周波数特性補正係数と周波数軸変換係数とに基づいて処理を行ない、マッチング部 60 は、このようにして周波数特性補正部 10、特徴量抽出部 20、周波数軸変換部 30 により得られた入力音声特徴量を標準音声記憶部 50 に記憶されている音声認識用の標準音声特徴量とマッチングして、最大尤度を与える標準音声特徴量に対応した語を認識結果候補として出力するようになっている。

【0015】

【発明が解決しようとする課題】この方法は、上記のように話者に発声内容を指定して発声を行わせることにより話者の個人性特徴を正規化させる最適な変換係数を推定する話者適応フェーズと、話者適応フェーズにおいて推定された変換係数を用いて未知内容の発声を認識する音声認識フェーズからなり、2つのフェーズを切り替えて使用するように構成されている。

【0016】しかしながら、「教師あり条件」での話者正規化手法であるこの従来法は、事前に未知話者に対し発声語彙を指定し学習データとして収録する必要があるため利用者への負担増を招いている。そこで、この負担を取り除くために、話者に事前に発声を要求しない「教師なし」条件で、かつ即効性のある方法での話者正規化方法の確立が必要である。

【0017】本発明では、発声者の発声器官の特徴つまり個人性を正規化することにより、ここでは声道長差に起因する影響を除去するために周波数軸変換を用いた「教師なし」条件での話者正規化を行う手法を述べる。具体的には、未知話者による内容未知の発声を用いて、周波数軸変換係数の精度よい推定方法を確立することで、教師なし条件の話者正規化方法を実現するものである。以下手法を実現するために生じる課題について述べる。

【0018】従来法などでは、声道長の違いは音声スペクトルが伸縮する形で現れることに着目して、入力音声スペクトルの周波数軸を変換し、音響モデルなどの標準となる音声スペクトルとの差を吸収する周波数軸変換を用いた話者正規化は効果があることを提示している。このとき、周波数軸変換の際に与える周波数軸変換係数が、声道長の違いに相当するともいえる。これを声道長正規化と呼ぶ。

【0019】一方、音声は、声道の形を変えることによって音韻の特徴を作り出されているため、同一話者であっても、発声される音韻によって声道長は異なっており、声道長推定値も一定の範囲で変化することが知られている「音声認識における個人差の学習法について、古井、日本音響学会音声研究会資料 S75-25、1975-11」。つまり、発声される音韻によって、異なる話者ではもちろんのこと同一話者内でも変動するものと考えられる。

【0020】したがって、声道長差に起因する音声スベ

クトルへの影響を除去するために、入力音声の周波数軸変換を行うにあたり、最適な周波数軸変換係数も、音韻によって変動していると考えられる。

【0021】そのため、最適な周波数軸変換係数を推定し次回の発声に利用できるように「教師なし」条件での話者正規化を考えた場合、今回発声された音韻と次回発声される未知の音韻の違いによって、今回推定された最適な周波数軸変換係数が必ずしも次回の発声には最適とはいえず、このことが未知の発声内容に対応して正規化を行うことを難しくする要因となっている。したがって、個人性の特徴の一つである声道長正規化を行うにあたって、未知話者による発声内容が未知の入力音声を教師信号として用いて オンライン「教師なし」話者正規化を行うためには、周波数軸変換係数の推定精度をより高める推定方法の確立が必要となる。

【0022】また、この従来法では、音声スペクトルに対して周波数軸を伸縮することによって、声道特性に関する個人性正規化を行うにあたって、この際入力音声区間全体に一律の変換係数を用いている。そのため、声道特徴に無関係な無声音の区間に対しても周波数軸変換を行うと、特徴量としての性質を失う原因になりかねず、認識結果に悪影響を及ぼすことも考えられる。そこで、入力音声区間全体に一律の変換係数を用いるのではなく、有声音区間に限って周波数軸変換を行うことにより、精度よく周波数軸変換係数の推定を行うことができると考えられる。

【0023】本発明の目的は、上記の問題点を解決し、利用者に対して 予め発声を要求せず、発声ごとに話者の音声から個人性を精度よく推定することにより、未知話者による発声内容が未知の入力音声にすばやく話者正規化できる話者正規化方法を備えた音声認識装置を提供することである。

【0024】

【発明が解決するための手段】本発明による音声認識方法および音声認識装置は、声道長差に起因するスペクトルの伸縮の影響を除去するため、入力音声のスペクトルに対して周波数軸変換を行なうことによる話者正規化方式を用いる。未知話者による発声内容が未知の入力音声を教師信号として、最尤推定により最適な周波数軸変換係数を決定した上で、音韻の違いによる周波数軸変換係数のばらつきを吸収するために、過去の周波数軸最適変換係数との平滑化を行うものである。さらに推定した周波数軸変換係数と過去の周波数軸最適変換係数とを比較することによって、話者が交代した場合とみなして平滑化を行うこともできる。

【0025】同一話者内での声道長の変動には限界があるため、未知話者による発声内容が未知の入力音声を教師信号として推定された周波数軸変換係数は、ばらつきはあっても一定の範囲内に収束することが期待できる。しかしこのとき、次のような問題が考えられる。(1)

周波数軸変換係数を変化させた場合に、高い尤度をかせぐ周波数軸変換係数の範囲と、マッチング処理によって出力される認識結果候補が発声内容と合致するつまり認識正解する周波数軸変換係数の範囲とは必ずしも一致しない。(2) 発声内容によって高い尤度をかせぐ周波数軸変換係数の範囲の分布が異なる。(3) マッチング処理によって出力される認識結果候補が発声内容と異なっている場合つまり誤認識した場合には、誤った内容に対して推定を行うことになってしまうため、マッチング処理によって出力される認識結果候補が発声内容として正しい場合つまり認識正解した場合に比べて、推定される周波数軸変換係数が異なる値になることがありうる。

【0026】そこで、過去の発声より求めた周波数軸変換係数との平均して平滑化することにより、推定された周波数軸変換係数のばらつきを吸収して、現在の話者への最適周波数軸変換係数が求められるものとする。

【0027】さらに、話者が交代した場合には、前回の発声から推定した最適周波数軸変換係数と今回の発声から推定した最適周波数軸変換係数との差が大きくなることが考えられる。このことを利用して話者が交代した場合には、平滑化処理を初期化するなどのこれに対処を行うことが可能である。

【0028】また、最適な周波数軸変換係数推定時には、音声区間検出手段から出力された無声音／有声音区間情報に同期して、周波数軸変換を行う。このため周波数軸変換係数推定にとって有効な音声区間あるいは音韻にのみ周波数軸変換を行うことから、精度よく周波数軸変換係数を推定できる。

【0029】以上より、発声内容によらず、事前の発声を必要としないオンライン「教師なし」話者正規化方法を実現することが可能な高性能な音声認識装置を提供することができる。

【0030】

【発明の実施の形態】本発明の請求項1に記載の発明は、入力音声の特徴量を抽出する特徴量抽出ステップと、前記入力音声の特徴量の周波数軸を少なくとも1つの周波数軸変換係数から構成される周波数軸変換係数列を用いて変換する周波数軸変換ステップと、前記周波数軸変換を施した入力音声の特徴量と予め複数の話者から学習した音響モデル特徴量とをマッチングし、認識結果候補を出力するマッチングステップと、前記認識結果候補のうち少なくとも1つから表現される音素系列に対して少なくとも1つの周波数軸変換係数から構成される最適な周波数軸変換係数列を推定する最適変換係数推定ステップと、前記求められた最適な周波数軸変換係数列と保持された過去に求められた周波数軸変換係数列とを平滑化し、新たな周波数軸変換係数列を更新・保持する変換係数平滑化ステップとを有するもので、発声者の音声特徴量から周波数軸上に現れる個人性を吸収することにより、認識率の向上させる作用を有する。

【0031】請求項2に記載の発明は、変換係数平滑化ステップは、少なくとも1つの周波数軸変換係数から構成される最新の周波数軸変換係数列と、保持された過去に求められた少なくとも1つの周波数軸変換係数から構成される周波数軸変換係数列とを比較することによって話者の交代を検知することを特徴とするもので、話者交代を検出した際には交代前話者の推定周波数軸変換係数の影響を受けないように周波数軸変換係数を初期化するなどによって、交代後話者への最適な周波数軸変換係数を新たに求め、話者間の周波数軸変換係数の差異による、認識率の低下を防ぐ作用を有する。

【0032】請求項3に記載の発明は、入力音声の特徴量を抽出する特徴量抽出ステップと、前記入力音声から無声音／有声音区間を弁別検出する音声区間検出ステップと、前記入力音声の特徴量の周波数軸を前記無声音／有声音区間情報に応じて周波数軸変換係数列を用いて変換する周波数軸変換ステップと、前記周波数軸変換を施した入力音声の特徴量と予め複数の話者から学習した音響モデル特徴量とをマッチングし、認識結果候補を出力するマッチングステップと、前記認識結果候補のうち少なくとも1つから表現される音素系列に対して少なくとも1つの周波数軸変換係数から構成される最適な周波数軸変換係数列を推定する最適変換係数推定ステップと、前記求められた最適な周波数軸変換係数列と保持された過去に求められた周波数軸変換係数列とを平滑化し、新たな周波数軸変換係数列を更新・保持する変換係数平滑化ステップとを有することを特徴とするもので、発声内容によって変動する推定変換係数のばらつきを抑えることにより、周波数軸変換を用いた話者正規化のより高い効果を与える作用を有する。

【0033】請求項4に記載の発明は、最適変換係数推定ステップは、前記認識結果候補を求めた際に使用した少なくとも1つの周波数軸変換係数から構成される周波数軸変換係数列を基に複数の周波数軸変換係数列候補を設定し、それぞれの周波数軸変換係数毎に、前記マッチングステップの認識結果候補のうち少なくとも1つから表現される音素系列に対して、周波数軸変換を施して得られた入力音声特徴量と予め複数の話者から学習した音響モデル特徴量とマッチングにより尤度を求め、求めた尤度のうちで最大尤度を与える少なくとも1つの周波数軸変換係数から構成される周波数軸変換係数列を選択することを特徴とするもので、事前に発声を行い話者の個人性を学習する適応フェイズなどを設けず、認識時の発声そのものから学習を行う「教師なし」話者正規化を実現する作用を有する。

【0034】請求項5に記載の発明は、変換係数平滑化ステップは、未知話者に最適な周波数軸変換係数を平滑化する際に、保持された過去に求められた周波数軸変換係数の平均値と最新の周波数軸変換係数との距離を求め、所定の距離を満たす場合にのみ最新の係数として採

用し、新たな周波数軸変換係数として出力することを特徴とするもので、発声内容によって変動する推定変換係数のばらつきを抑えることにより、周波数軸変換を用いた話者正規化のより高い効果を与える作用を有する。

【0035】請求項6に記載の発明は、変換係数平滑化ステップは、未知話者に最適な周波数軸変換係数を平滑化する際に、保持された過去に求められた周波数軸変換係数の平均値と最新の周波数軸変換係数との距離を求め、所定の距離を満たさない場合には、過去の周波数軸変換係数を任意の初期値を与えることによって初期化を行うことを特徴とするもので、発声内容によって変動する推定変換係数のばらつきを抑えることにより、周波数軸変換を用いた話者正規化のより高い効果を与える作用を有する。

【0036】請求項7に記載の発明は、入力音声の特徴量を抽出する特徴量抽出手段と、前記入力音声の特徴量の周波数軸を周波数軸変換係数列を用いて変換する周波数軸変換手段と、前記周波数軸変換を施した入力音声の特徴量と予め複数の話者から学習した音響モデル特徴量とをマッチングし、認識結果候補を出力するマッチング手段と、前記認識結果候補のうち少なくとも1つから表現される音素系列に対して、少なくとも1つの周波数軸変換係数から構成される最適な周波数軸変換係数列を推定する最適変換係数推定手段と、前記求められた最適な周波数軸変換係数列と保持された過去に求められた周波数軸変換係数列とを平滑化し、新たな周波数軸変換係数列を更新・保持する変換係数平滑化手段とを有することを特徴とするもので、話者の音声特徴量から周波数軸上に現れる個性を吸収する話者正規化によって、認識率の向上させる作用を有する。

【0037】請求項8に記載の発明は、変換係数平滑化手段は、少なくとも1つの周波数軸変換係数から構成される最新の周波数軸変換係数列と、保持された過去に求められた少なくとも1つの周波数軸変換係数から構成される周波数軸変換係数列とを比較することによって話者の交代を検知することを特徴とするもので、話者交代を検出した際には交代前話者の推定周波数変換係数の影響を受けないように周波数軸変換係数を初期化するなどによって、交代後話者への最適な周波数軸変換係数を新たに求め、話者間の周波数軸変換係数の差異による、認識率の低下を防ぐ作用を有する。

【0038】請求項9に記載の発明は、入力音声の特徴量を抽出する特徴量抽出手段と、前記入力音声から無声音／有声音区間を弁別検出する音声区間検出手段と、前記入力音声の特徴量の周波数軸を前記無声音／有声音区間情報に応じて周波数軸変換係数列を用いて変換する周波数軸変換手段と、前記周波数軸変換を施した入力音声の特徴量と予め複数の話者から学習した音響モデル特徴量とをマッチングし、認識結果候補を出力するマッチング手段と、前記認識結果候補のうち少なくとも1つから

表現される音素系列に対して、少なくとも1つの周波数軸変換係数から構成される最適な周波数軸変換係数列を推定する最適変換係数推定手段と、前記求められた最適な周波数軸変換係数列と保持された過去に求められた周波数軸変換係数列とを平滑化し、新たな周波数軸変換係数列を更新・保持する変換係数平滑化手段とを有することを特徴とするもので、発声内容によって変動する推定変換係数のばらつきを抑えることにより、周波数軸変換を用いた話者正規化のより高い効果を与える作用を有する。

【0039】以下、本発明の実施の形態について図を用いて説明する。

【0040】（実施の形態1）図1は、本発明の実施の形態1における音声認識装置のブロック図である。図1において、1は入力音声に対してA/D変換処理などを行う音声取り込み手段、2は音声の音響的特徴をモデル化した音響モデル、3は単語系列における単語間の関係をモデル化した言語モデル、4はデータやプログラム装置に入力する入力手段、5はデータやプログラムを記録するメモリ、6はプログラムにしたがってデータを処理したり装置全体を制御するCPU、7は認識結果候補を出力する出力手段である。

【0041】図2は、本発明の音声認識装置の処理手順を示すフローチャートであり、この図を用いて音声認識装置の処理手順を説明する。

【0042】入力音声の特徴量抽出が行われるS10では、マイクロフォン等から取りこまれた音声にA/D変換を行ってデジタル化された音声を入力信号として、一定フレーム周期毎にLPCメルケプストラム係数を出力する。具体的には、文献「音声認識、今井著、共立出版、1995年11月25日」などに示されているこの方法を利用して、プリエンファシス： $1-z^{-1}$ 、窓周期：20ms、フレーム周期：10ms、LPC分析次数：10次、ケプストラム分析次数：10次としてLPCメルケプストラム係数を出力する。入力音声はここでは、8kHzサンプリングされるものとする。

【0043】音声区間検出処理が行われるS20では、無声音／有声音区間などの検出を行なう。たとえば、入力音声信号のフレームパワーを求め、入力開始数フレームでの平均値を求めておき、その平均値に比べ、フレームパワーが2倍になったフレーム区間を有声音区間とする方法などを用いることにより、有声音区間とすることもできる。あるいは、入力音声に1500Hz-3400Hz通過高域フィルタを掛けて、同様に高域音声入力があったことを示す高域区間を求めておき、これを利用することにより、有声音区間、無声音区間の判定を大まかな判定をすることができる。

【0044】周波数軸変換処理が行われるS30では、声道長の個人差に起因するスペクトルを伸縮することによって声道長正規化を行う周波数軸変換を施す。具体的

には、入力音声のスペクトルを表現している L P C メルケプストラム係数に対して、例えば (数 1) で表わされる 1 次の全域透過フィルタ $H(z)$ を作用させて、周波数軸の変換を行なう。この周波数軸変換の手法は、論文「LPC 距離尺度における周波数正規化に関する検討、小

$$H(z) = (z^{-1} - \alpha) / (1 - \alpha z^{-1}) \quad (-1 < \alpha < 1)$$

【0046】また、実際に (数 1) を用いて周波数軸を変換によって、スペクトル伸縮による周波数軸変換後の L P C メルケプストラム係数の算出方法としては、たとえば、論文「Discrete representation of signals, Oppenheim and Johnson, Proc. IEEE, 60, pp681-691, June 1972」で示されている手法を用いる。

【0047】なお、本実施例では、L P C メルケプストラムの次数を例えば 10 次としたメル周波数変換の処理も同時に行なう。メル尺度を最も良く近似する周波数軸変換係数 α の値は、サンプリング周波数 8 kHz の場合、 $\alpha_0 = 0.315$ 付近とされており、この値を α の値を基準として $\alpha = \alpha_0 + 0.05$ 、 $\alpha = \alpha_0 - 0.05$ などと指定してスペクトルを伸縮させる。ここで指定する周波数軸変換係数の値は、変換係数平滑化処理が行われる S 70 から算出された値を用いる。

【0048】このとき、周波数軸変換係数の推定によって有効な区間あるいは音韻のみ周波数軸変換を行うことによって、精度よく変換係数を推定させるため、音声区間検出処理から出力された無声音／有声音区間情報に同期して周波数軸変換を行う。たとえば、音声区間検出処理の行われる S 20 より得た有声音区間のフレームにのみ、周波数軸変換を実施する。

【0049】マッチング処理が行われる S 50 では、S 30 において周波数軸変換された入力音声特徴量と予め複数の話者から学習した音響モデル特徴量とを、言語モデルとして表現される単語辞書等を用いてマッチングを行い、少なくとも 1 つの認識結果候補を出力する。なお、このマッチングは、例えば、端点フリー DP (ダイナミック・プログラミング) マッチング法によりなされる。

【0050】最適変換係数推定処理が行われる S 60 では、マッチングにより出力される認識結果候補のうち少なくとも 1 つから表現される音素系列に対して、複数の周波数軸変換係数候補を設けて、それぞれの周波数軸変換係数毎に、周波数軸変換された入力音声の特徴量と予め複数の話者から学習した音響モデル特徴量とのマッチングを行い尤度を求めて、最尤推定により最大尤度を与える周波数軸変換係数を最適な周波数軸変換係数 α_0 として決定し、平滑バッファに登録する。

【0051】この時、最適変換係数推定処理が行われる S 60 で推定される最適な周波数軸変換係数は発声内容によってばらつきが生じたり、不正解の音素系列に対して推定を行った場合には、必ずしも最適でない場合もありうるということが問題となる。

林 松本 熊田、1983、日本音響学会音声研究会資料 S83-47, 1983 Dec. 22」に示されているものである。

【0045】

【数 1】

【0052】そこで、変換係数平滑化処理が行われる S 70 では、推定した最適な周波数軸変換係数のばらつきを吸収するために、最適変換係数推定処理 S 60 で求められた周波数軸変換係数 α_0 と過去の周波数軸変換係数を記憶した平滑化バッファより読み出された、たとえば過去 10 回の周波数軸変換係数の平均により平滑化された周波数軸変換係数を算出し、新たな周波数軸変換係数 α_0 として更新・記憶する。

【0053】次に、最適変換係数推定処理について、図 3 (a) (b) を用いて詳細に説明する。図 3 (a) で示されるように、変換係数候補の値として、認識結果候補を求めた際用いた周波数軸変換係数 α_0 に対して、 $\alpha_0 - 0.05$ 、 α_0 、 $\alpha_0 + 0.05$ の 3 点を設定し (S 101)、それぞれ周波数軸変換を施した入力音声特徴量と認識結果候補第一位 $r(1)$ の単語の音素系列で表現される音響モデル特徴量系列とから尤度を求める (S 102)。たとえば、周波数軸変換係数候補を、 $x_0 = \alpha_0 - 0.05$ 、 $x_1 = \alpha_0$ 、 $x_2 = \alpha_0 + 0.05$ を設定し、それぞれ得られる尤度を y_0 、 y_1 、 y_2 とする。

【0054】

【数 2】

$$z_1 = (y_1 - y_0) / (x_1 - x_0)$$

$$z_2 = (y_2 - y_0) / (x_2 - x_0)$$

$$a_k = (z_2 - z_1) / (x_2 - x_1)$$

$$b_k = z_1 - (x_0 + x_1) * a_k$$

【0055】周波数変換係数およびその尤度からなる 3 点を用いて二次曲線に近似すると、その二次近似曲線がピークをとる周波数変換係数 α_x は、次のようにあらわされる (S 103)。

【0056】

【数 3】

$$\alpha_x = -b_k / 2/a_k$$

【0057】さらに、前記 3 点同様に周波数軸変換係数 α_x に対する尤度を求め (S 104)、4 点の中から最大尤度を与える周波数軸変換係数 α_0 を採用する (S 105) のものである。

【0058】次に、尤度距離の計算処理を図 3 (b) の処理フローチャートを用いて説明する。

【0059】変換係数候補の値として設定された周波数軸変換係数に対して、入力音声の特徴量に周波数軸変換を施す (S 106)。前記の周波数軸変換を施した入力音声の特徴量と、認識結果候補第 1 位 $r(1)$ の単語の音素

系列から表現される音響モデル特徴量の系列とから尤度を求める (S107)。

【0060】なお、これまで最適変換係数推定の説明において、認識結果候補のうち第一位の候補 $r(1)$ のみを用いたが、これを認識結果候補 $r(n)$ を n 位まで利用して、最大尤度を与える周波数軸変換係数を採用することもできる。

【0061】また、上記の特徴量を周波数軸変換する際に、音声区間検出処理が行われるS20で算出される無声音/有声音などの音声区間情報に同期して、複数の周波数軸変換係数を用いることもできる。たとえば、無声音の区間にのみ周波数軸変換を適応したり、逆に有声音の区間にのみ周波数軸変換を適応したり、音声パワーの同じ区間ごとに異なる周波数軸変換係数を用いることもできる。

【0062】なお、変換係数平滑化処理において、最適変換係数推定処理S60で算出された周波数軸変換係数を平滑化バッファに登録する際に、発声される音韻による声道長の変動に伴う最適な周波数軸変換係数のばらつきを吸収するため、さらには推定精度の低下を防ぐため、今回の発声に対して推定された最適な周波数軸変換係数を評価し、平滑化バッファに登録するか否かを判断する。その詳細について、図4を用いて説明する。

【0063】前記の認識結果候補を求めた際用いた周波数軸変換係数 α_0 つまり前回までの発声に対して推定された最適な周波数軸変換係数 α_0 と、今回の発声に対して推定された最適な周波数軸変換係数 α_n を比較するにあたって、たとえば、 $|\alpha_0 - \alpha_n| < 0.100$ という評価関数 (S110、S111) を用いて、この条件を満たす場合にのみ今回の発声に対して推定された周波数軸変換係数 α_n を採用・登録し、平滑化バッファ内の周波数軸変換係数を平均化し、周波数軸変換係数 α_0 を更新する (S112)。この周波数軸変換係数 α_0 は、次の発声に対して、図2のS30などでの周波数軸変換に用いられることになる。

【0064】このように評価関数を用いることにより、推定された周波数軸変換係数のばらつきを抑える事ができることから、平滑化された周波数軸変換係数を精度よく求めることができるため、周波数軸変換を用いた話者正規化による効果をより高めることが可能となる。

【0065】(実施の形態2) 上記実施の形態1の音声認識装置は、一人の話者が発声していることを前提としているため、現在の発声話者に対して最適な周波数軸変換係数が推定され、この変換係数を用いて次の発声に対して話者正規化が実施される。一方、家庭など複数の話者が交代で利用することが想定される場でこの音声認識装置を用いると、交代前話者に対して推定された周波数軸変換係数を用いて、交代後の話者に対して話者正規化が行われることになる。この時、交代直後に推定される周波数軸変換係数は 交代後の話者にとって必ずしも最

適な値とは限らず、認識率を低下させる原因にもなり兼ねない。

【0066】そこで、本発明の実施の形態2では、このように複数の話者が交代するような状況においては、話者交代を検出することによって、交代前話者の推定周波数軸変換係数の影響を受けずに、交代後話者に対して最適な周波数軸変換係数を求めることを可能とするものである。

【0067】本実施の形態の音声認識装置における話者交代を検知する実施形態について、図6のフローチャートを用いて説明する。本実施形態は、実施の形態1の図2で説明した実施例に加えて、話者交代を検知する手段を有することを特徴としているもので、実施の形態1と説明の重複を省くために異なる部分のみを説明する。

【0068】図2同様、S60にて今回の発声に対して推定された最適な周波数軸変換係数に対して、過去所定回数の発声に対して推定された最適な周波数軸変換係数から算出される現話者に対して推定された最適な周波数軸変換係数とを比較して、今回の周波数軸変換係数を評価することにより話者の交代を検知し、話者交代を検知した場合には、前回までの発声に対して推定された最適な周波数軸変換係数保持している平滑化バッファを初期化し、新しい話者の周波数軸変換係数の登録を行う (S80)。最後の処理として、図2同様に、現話者に対して推定された最適な周波数軸変換係数を算出する (S70)。

【0069】ここで、S80における話者交代検出話者の処理の詳細について、図5を用いて説明する。話者の交代の検出は、前回の発声に対して推定された周波数軸変換係数 α_{n-1} と今回の発声に対して推定された周波数軸変換係数 α_n を比較して行い、たとえば $|\alpha_{n-1} - \alpha_n| < 0.150$ という評価関数 (S120、S121) を用いる。この条件を満たさない場合には、話者の交代を検出したとみなし、平滑化バッファの初期化を行うことにより (S122)、新しい話者に対する現話者に対して推定された最適な周波数軸変換係数 α_n を出力する。

【0070】このようにして、話者が交代しても認識率を低下させず、しかも従来法のように指定された単語の発声を促すこともなく、話者正規化を進めることが可能となる。

【0071】(実施の形態3) 次に、本発明の実施の形態3の音声認識装置について説明する。

【0072】実施の形態1または実施の形態2では、一段階のマッチング法を用いたが、実施の形態3では2段階の認識方式について説明する。

【0073】本発明の実施の形態3における音声認識装置の動作フローチャートを図7に示すが、実施の形態1または実施の形態2と異なる部分についてのみ説明するものとする。

【0074】1段目の予備マッチングを行い (S50)、

出力される認識結果候補のうち少なくとも1つから表現される音素系列に対して最適な周波数軸変換係数を推定した後(S60)、話者交代を考慮して(S80)、現話者に対して推定された最適な周波数軸変換係数を算出する(S70)。さらにS70にて求めた最適な周波数軸変換係数を用いて周波数軸変換を施した(S31)特徴量を用いて、2段目の精密マッチング(S51)を行う実施形態も可能である。

【0075】(実施例)以上、本実施例の構成を用いて、100単語を発声した男女50名の音声データの認識実験を行った。まず、オンライン「教師なし」話者正規化を実現するために、変換係数平滑化の効果を調べる実験を行った。この実験においては、あらかじめ話者正規化を行わない条件で男女50名について認識率を算出して、その結果認識性能の悪い10話者を対象とした。

【0076】図8に示す実験結果より、話者正規化学習

は、評価発声データを7単語以上ではその効果がほぼ飽和していることから、評価発声データ10単語を単位とすれば、オンライン「教師なし」話者正規化には十分効果があることが分かる。

【0077】次に、音声区間情報に同期して周波数軸変換を行う効果について男女50名について調べる実験を行った。なお、変換係数平滑化のためのオンライン学習データ数は10とした。その結果(表1)に示すように、話者正規化を行う前は93.76%であったのに対し、音声区間情報を用いずに話者正規化を一律に行った場合、94.78%、本実施例に基づく有声音区間のみ話者正規化を行った場合は、95.44%に認識率が改善され、誤り率もそれぞれ約16.0%、約26.9%、改善された。

【0078】

【表1】

	50名平均	最低10名	最低話者
話者正規化なし	93.76%	86.9%	77%
話者正規化あり	94.78%	89.4%	83%
話者正規化あり[有声音区間のみ正規化]	95.44%	91.2%	87%

【0079】また、50名の中で認識率の悪い10名についても効果が認められ、最低話者についても77%から、83%、87%と大幅に改善された。このことから、音声区間情報に同期して周波数軸変換を行なうことによる効果が認められる。

【0080】なお、本実施例においては、単語マッチング方法として端点フリーのDPマッチング法を用いたが、HMM(隠れマルコフモデル)での実施も可能である。

【0081】なお、本実施例においては、単語マッチング時の距離尺度として共分散行列を共通化したマハラノビス距離を用いたが、共分散行列を共通化しないマハラノビス距離や、HMMから構成される音素モデルからを用いて計算することもできる。

【0082】また、本実施例においては、認識対象を単語としたが、これを連続発声認識する際に利用することも可能である。

【0083】なお、本実施例においては、音響的特徴を表現する特徴量としてLPCメルケプストラム係数を用いたが、LPCケプストラム係数、メルケプストラム係数、ケプストラム係数での実施も可能である。

【0084】なお、本実施例においては、入力される音声は、8kHzでサンプリングされたデータを用いたが、他のサンプリング周波数についても実施可能である。なお、本実施例においては、マッチングの際に音声区間情報は用いなかったが、音声区間情報を用いて、音声の始端を制限したマッチングなどを行うこともできる。

【0085】

【発明の効果】本発明によれば、声道長差に起因するスペクトルの伸縮の影響を除去するため、入力音声のスペ

クトルに対して周波数軸変換を行なうにあたり、過去所定回数の発声から推定した変換係数の平均値を用いて、推定された周波数軸変換係数のばらつきの影響を抑えることにより、オンライン「教師なし」条件で話者正規化を実現し、高性能な音声認識装置を提供できるという効果を得る。さらに、有声音/無声音などの音声区間情報を利用して、声道特性に無関係な区間に対しては周波数軸変換を行わず有声音区間のみ周波数軸変換を行うことにより、周波数軸変換係数の推定をより精度のよく行なうことができるため認識率の向上を図ることが可能である。

【図面の簡単な説明】

【図1】本発明の実施の形態における音声認識装置のブロック図

【図2】本発明の実施の形態1における音声認識装置の処理フローチャート

【図3】(a)最適変換係数推定処理のフローチャート
(b)尤度計算処理フローチャート

【図4】変換係数平滑化処理のフローチャート

【図5】話者交代検出処理のフローチャート

【図6】本発明の実施の形態2における音声認識装置の話者交代検出処理のフローチャート

【図7】本発明の実施の形態3における音声認識装置の2段階認識処理のフローチャート

【図8】変換係数学習データバッファの大きさと認識率を示す図

【図9】従来例の音声認識装置の処理フローチャート

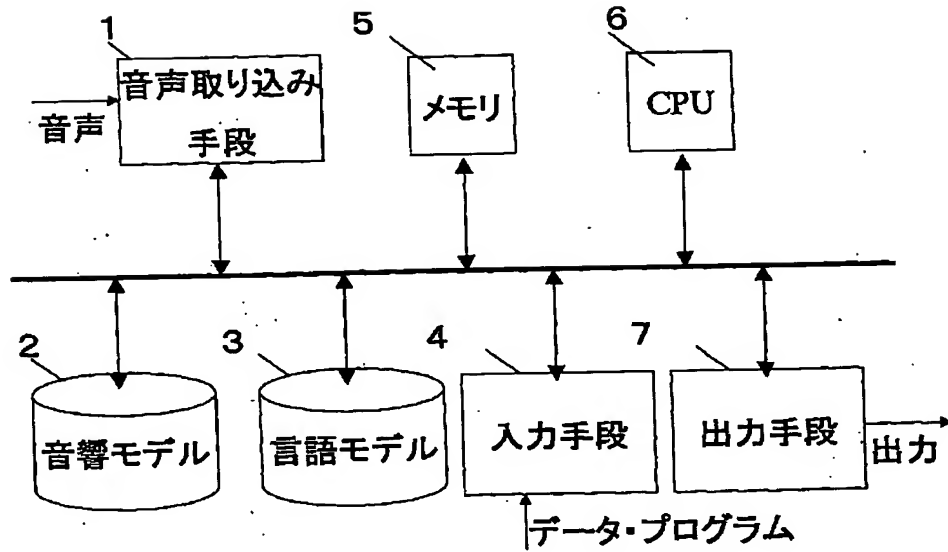
【符号の説明】

- 1 音声取り込み手段
- 2 音響モデル

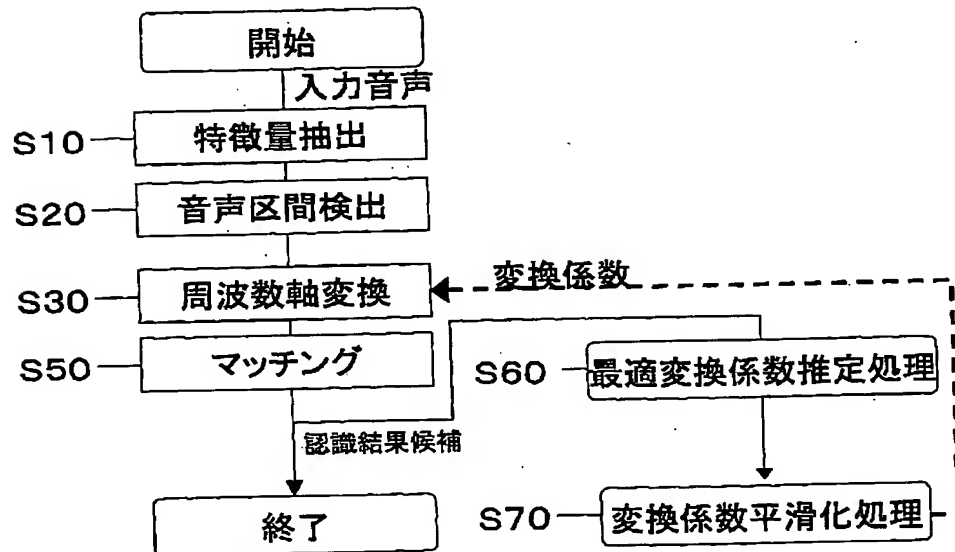
- 3 言語モデル
- 4 入力手段
- 5 メモリ

- 6 CPU
- 7 出力手段

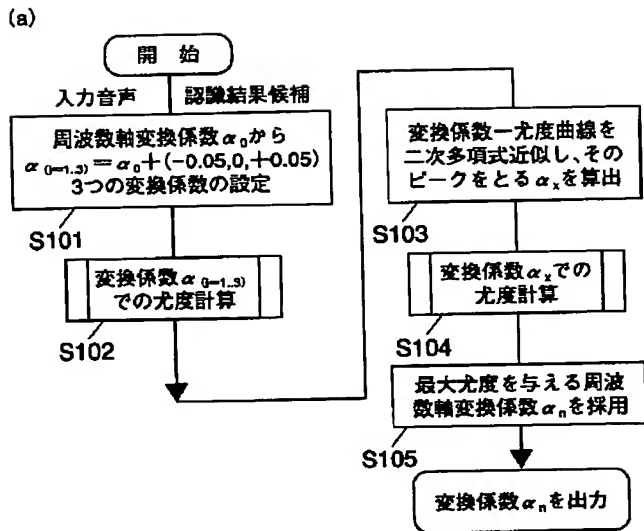
【図1】



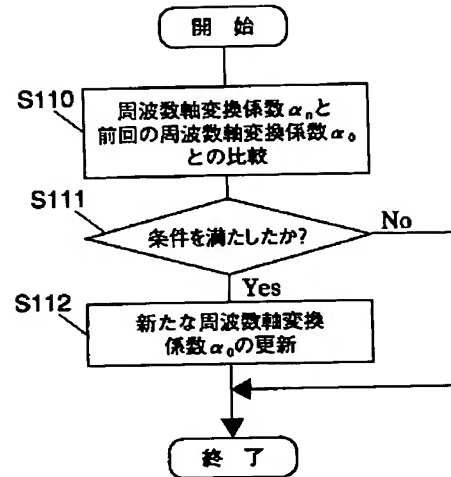
【図2】



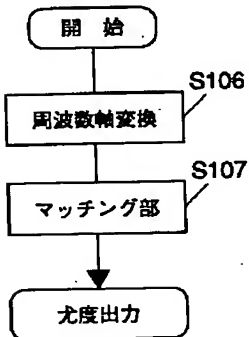
【図3】



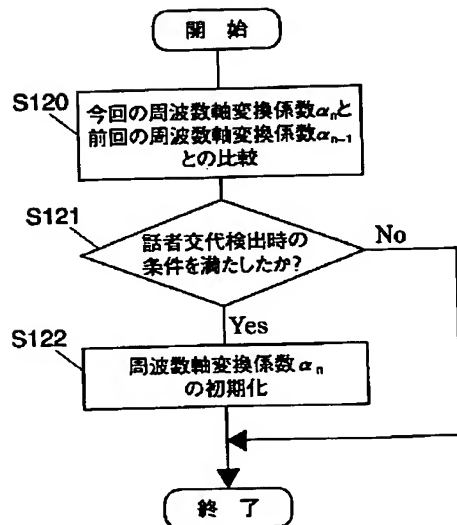
【図4】



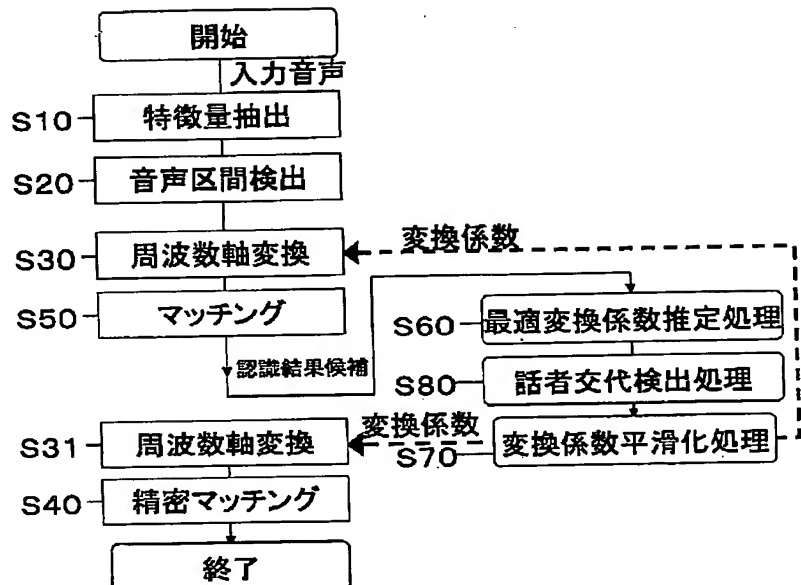
(b)



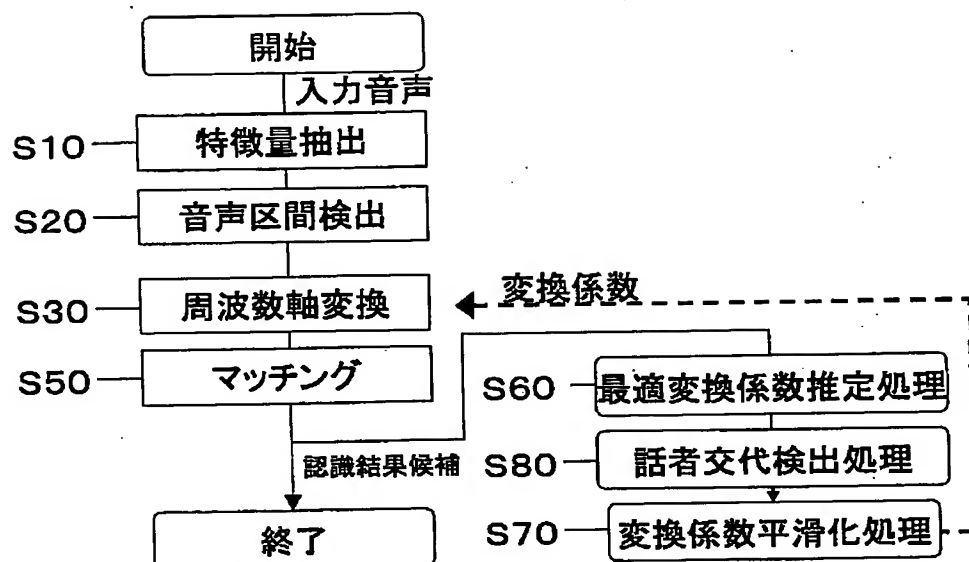
【図5】



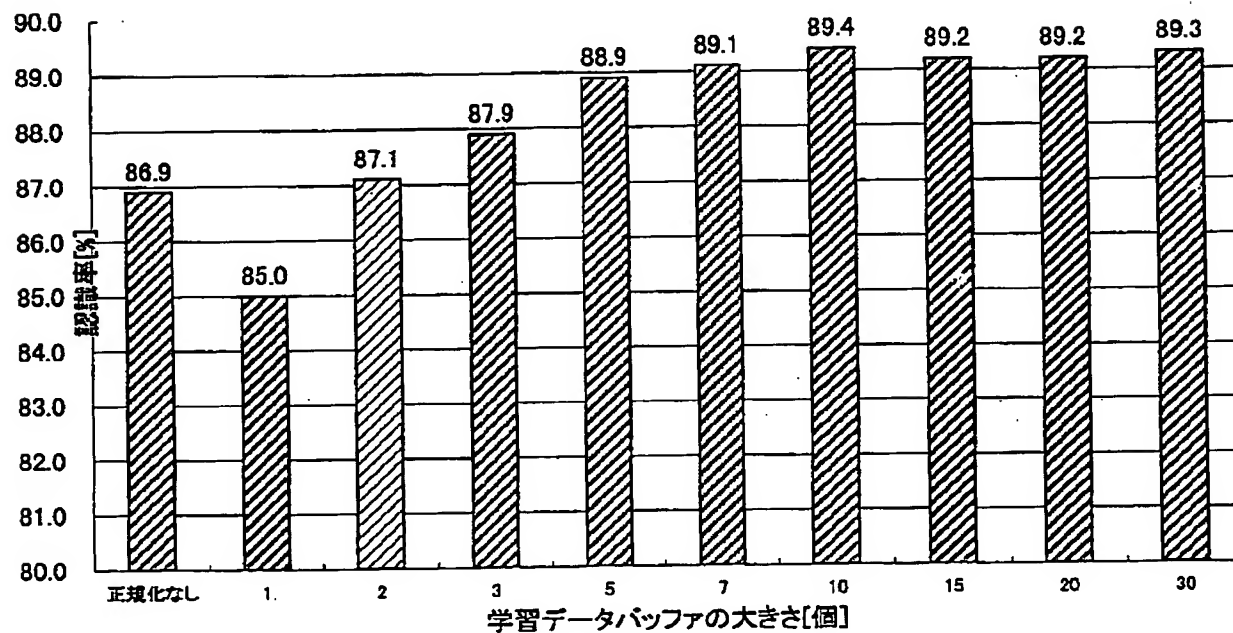
【図7】



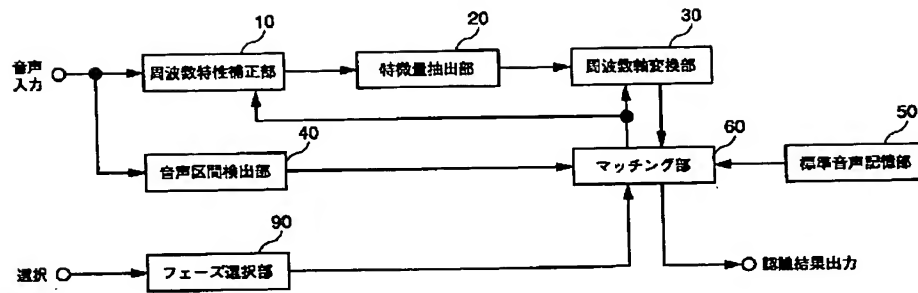
【図6】



【図8】



【図 9】



フロントページの続き

(51) Int. Cl.⁷
// G 1 0 L 101:16

識別記号

F I

テ-マコ-ト (参考)

(72) 発明者 木村 達也
神奈川県川崎市多摩区東三田 3 丁目 10 番 1
号 松下技研株式会社内

F タ-ム (参考) 5D015 AA02 BB02 FF07
9A001 BB06 EE05 GG01 HH16 HH17